

# 从 Transformer 到通用时间序列预测 From Transformer to Universal Forecaster

李丰

北京大学光华管理学院

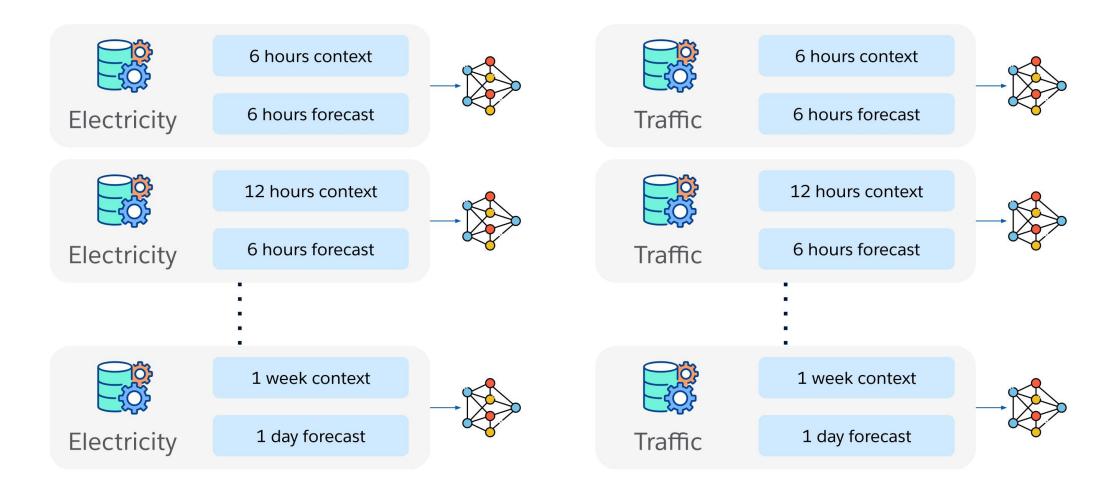
https://feng.li/forecasting-with-ai

### 传统的预测范式

Existing forecasting paradigm



#### 传统时间序列预测范式



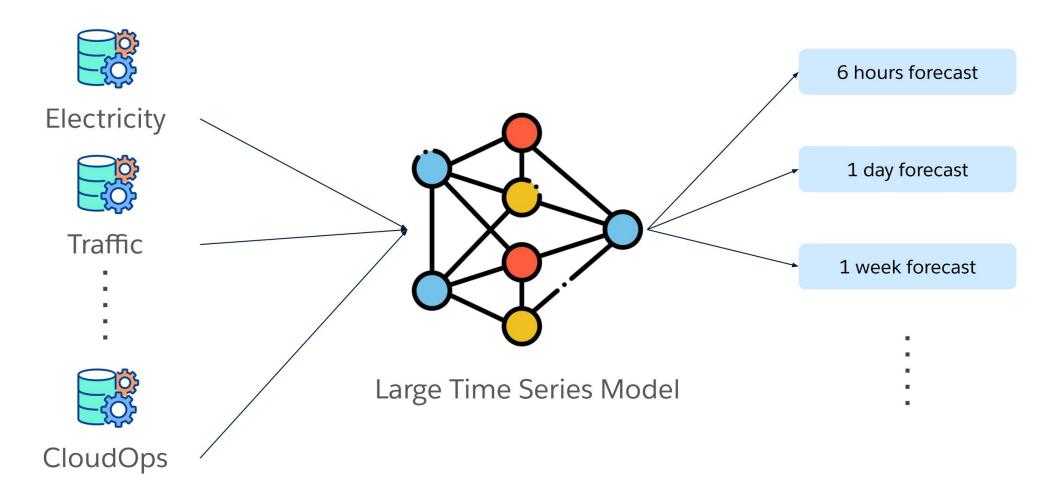


#### 局限与挑战

- 当前主流的深度预测方法通常遵循以下模式:
  - "一模型对应一数据集、一上下文长度、一预测长度"
  - 每个数据集都需要单独训练模型
  - 模型只能处理固定的输入窗口与固定的预测步长
  - 当数据集或任务设置发生变化时,模型需重新训练。
- 典型模型举例: ARIMA, DeepAR、N-BEATS、Temporal Fusion Transformer (TFT)、Informer
- 主要局限:
  - 可扩展性差(Scalability Issue)
  - 不同数据集、预测任务需重复训练, 计算成本高。
  - 泛化能力弱(Poor Generalization)
  - 模型难以跨领域、跨时间尺度迁移。
  - 灵活性不足(Low Flexibility)
  - 资源浪费(Resource Inefficiency)
- 每次任务变化都要重新建模,参数无法共享。



#### 全新预测范式: 通用预测模型 (Universal Forecaster)





#### 通用预测模型

- 预训练+微调 (Pretrain & Fine-tune)
- 支持多任务与可变上下文长度
- 代表模型:
  - PatchTST: Nie, Y., Nguyen, N. H., Sinthong, P., & Kalagnanam, J. (2023). A Time Series is Worth 64 Words: Long-term Forecasting with Transformers <a href="https://doi.org/10.48550/arXiv.2211.14730">https://doi.org/10.48550/arXiv.2211.14730</a>
  - TimeGPT: Garza, A., Challu, C., & Mergenthaler-Canseco, M. (2024). TimeGPT-1 <a href="https://doi.org/10.48550/arXiv.2310.03589">https://doi.org/10.48550/arXiv.2310.03589</a>
  - Chronos: Ansari, A. F., Stella, L., Turkmen, C., Zhang, X., Mercado, P., Shen, H., Shchur, O., Rangapuram, S. S., Arango, S. P., Kapoor, S., Zschiegner, J., Maddix, D. C., Wang, H., Mahoney, M. W., Torkkola, K., Wilson, A. G., Bohlke-Schneider, M., & Wang, Y. (2024). Chronos: Learning the Language of Time Series. <a href="https://doi.org/10.48550/arXiv.2403.07815">https://doi.org/10.48550/arXiv.2403.07815</a>
  - Moirai: Woo, G., Liu, C., Kumar, A., Xiong, C., Savarese, S., & Sahoo, D. (2024). Unified Training of Universal Time Series Forecasting Transformers. <a href="https://doi.org/10.48550/arXiv.2402.02592">https://doi.org/10.48550/arXiv.2402.02592</a>



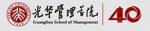
#### Transformer 在通用预测模型中的主导地位

- Transformer 自注意力机制(Self-Attention)能捕捉长程依赖
- 输入长度灵活,天然适配变长序列
- 可并行计算,训练效率高
- 结构通用性强,语言、图像、时间序列均可适用。

模型	核心结构	特点
PatchTST (2023)	Vision Transformer (ViT) 思想	将时间序列切片为"时间块"作为Token
TimeGPT (Nixtla, 2024)	Decoder-only Transformer	基于大规模多领域时间序列预训练
Chronos (Amazon, 2024)	GPT架构 + Token化时间序列	使用离散化时间Token和概率建模
Moirai (Salesforce, 2024)	Sparse Transformer + Mixture-of-Experts (MoE)	支持可变步长和多频率输入
TimesFM (Google, 2024)	Encoder-only Transformer	在数百万序列上预训练

### Transformer

Attention Is All You Need



#### 从序列出发:语言与时间的共同逻辑

- 企业数据中处处是"序列":
  - 客户行为日志(点击、购买、退货)
  - 财务指标(季度报表)
  - 市场舆情(时间演化)
- 传统方法(ARIMA, RNN、LSTM)的问题:
  - 无法并行 (scale law)
  - 记忆短期,忽视长期关系
- Transformer 突破点: 同时关注"整个序列"
- 示例: 预测顾客是否复购,既要看"最新点击"也要看"历史购买模式"。

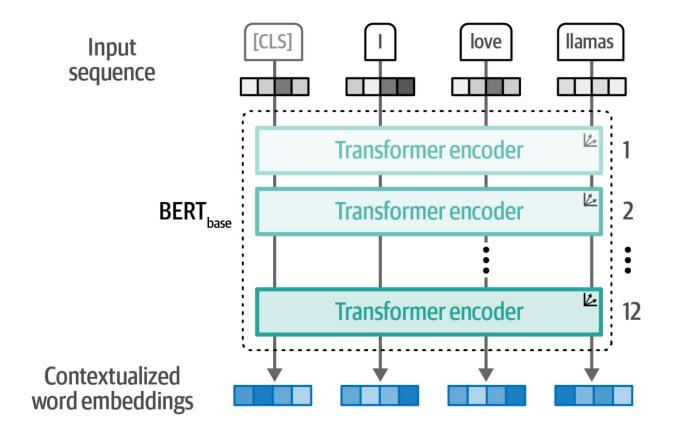


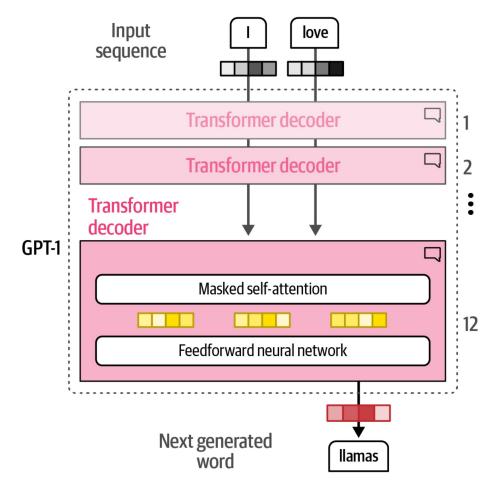
#### Transformer 的注意力机制

- 概念: 注意力机制(Attention)= 让模型学会"关注重点"  $Attention(Q,K,V) = softmax(QK^T/\sqrt{d_k})V$
- *Query* (问题): 你在问"谁最重要?"
- *Key*(线索):每个人说了什么
- *Value* (答案): 对应的信息内容
- 模型根据相关性决定"谁值得被听"
  - $QK^T$ : 计算相似度: 每个 Query 与所有 Key 相乘得到"相似程度"。值越大表示 Query 对该 Key 越"关注"。
  - 除以 $\sqrt{d_k}$ :缩放避免维度过大导致值过大,致使 softmax 输出过于极端(梯度不稳定)。
  - softmax: 归一化注意力权重,把所有相似度转成概率权重(总和为 1)
  - 乘以 V: 加权求和得到注意力输出。每个输出是所有 Value 的加权组合,权重来自 attention 分数。



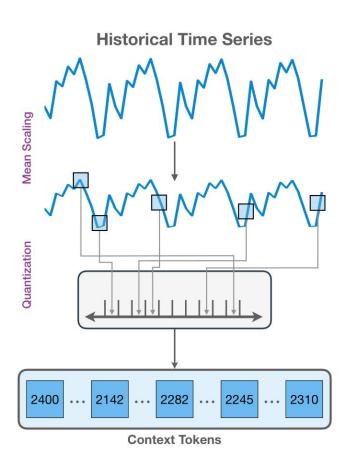
#### 从 BERT 到第一代 GPT 模型

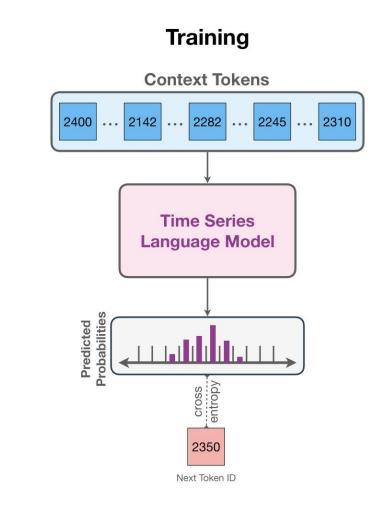




#### 学习时间序列的语言

#### **Time Series Tokenization**





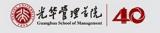
### Inference **Context Tokens** 2142 2282 ... | 2245 | ... | 2310 **Time Series Language Model** 2283 Dequantization and Unscaling

**Probabilistic Forecast** 

#### 通用时间序列预测模型对比

模型	任意变量(零样本)	是否支持概率预测	分布灵活性	预训练数据(规模)	是否开源
MOIRAI	√ 是	√ 是	√ 是	LOTSA(> 270亿条)	√ 是
TimeGPT-1	√ 是	√ 是	<b>X</b> 否	未公开(约1000亿条)	<b>X</b> 否
ForecastPFN	<b>X</b> 否	<b>X</b> 否	-	合成数据(6000万条)	√ 是
Lag-Llama	<b>X</b> 否	√ 是	<b>X</b> 否	Monash 数据集(< 10亿 条)	√ 是
TimesFM	<b>X</b> 否	<b>X</b> 否	-	Wiki + Trends + 其他数据 (> 1000亿条)	√ 是
TTM	<b>X</b> 否	<b>X</b> 否	-	Monash 数据集(< 10亿 条)	√ 是
LLMTime	<b>X</b> 否	√ 是	√ 是	网络级文本(Web-scale Text)	√ 是

- 这些通用时间序列模型的预训练数据规模差异巨大,从数千万到千亿级不等。
- 源涵盖合成数据、学术基准集、网络文本与多源趋势数据。
- 总体来看,数据多样性与规模是决定模型泛化能力与零样本预测性能的关键因素。



## 预测模型的性能比较



#### TimeGPT-1 与传统模型的预测性能比较

	Monthly rMAE rRMSE		Weekly rMAE rRMSE		Daily   rMAE rRMSE		Hourly   rMAE rRMSE	
ZeroModel	2.045	1.568	6.075	6.075	2.989	2.395	10.255	8.183
HistoricAverage SeasonalNaive	1.349 1.000	1.106 1.000	4.188 1.000	4.188 1.000	2.509 1.000	2.057 1.000	2.216 1.000	1.964 1.000
Theta	0.839	0.764	1.061	1.061	0.841	0.811	1.163	1.175
DOTheta ETS	0.799 0.942	0.734 0.960	1.056 1.079	1.056 1.079	0.837 0.944	0.806 0.970	1.157 0.998	1.169 1.009
CES	1.024	0.946	1.002	1.002	0.919	0.899	0.878	0.896
ADIDA IMAPA	0.852	0.769 0.769	1.364	1.364 1.364	0.908	$0.868 \\ 0.868$	2.307 2.307	2.207 2.207
CrostonClassic	0.989	0.857	1.805	1.805	0.995	0.933	2.157	2.043
LGBM	1.050	0.913	0.993	0.993	2.506	2.054	0.733	0.709
LSTM	0.836	0.778	1.002	1.002	0.852	0.832	0.974	0.955
DeepAR TFT	0.988	$0.878 \\ 0.700$	0.987	0.987 0.954	0.853	0.826 0.791	1.028 1.120	1.028 1.112
NHITS	0.738	0.694	0.883	0.883	0.788	0.771	0.829	0.860
TimeGPT	0.727	0.685	0.878	0.878	0.804	0.780	0.852	0.878

- 在所有频率上,TimeGPT 的 rMAE 与 rRMSE 均为最小值或 接近最小值。
- 传统统计模型在低频数据上仍 具竞争力。如 Theta 和 DOTheta 在月度、周度任务中 表现优于其他传统方法。ETS、 CES 在日度任务中也能保持较 好性能。
- · · 深度模型内部差异显著。在深 · · · 度模型中,NHITS 与 TFT 表现 明显优于 LSTM 与 DeepAR。
  - 整体趋势频率越高,误差越大。 反映高频序列的噪声性和建模 难度更高。



#### MOIRAI 与传统模型的预测性能比较

- 多种时间序列预测模型在不同数据集下的概率预测表现。
- 概率预测评估指标为 CRPS(连续分级概率评分,越低越好)与 MSIS(平均尺度化区间得分,越低越好)。
- 结果区分了 Zero-shot (零样本预测)、Full-shot (完全训练)与 Baseline (基线模型)三类。
- 主要结论
  - MOIRAI 系列模型在 Zero-shot 预测中表现最优
  - 传统方法在概率预测上已明显落后
  - 预训练+注意力架构正成为时间序列预测的新范式

		<b>Zero-shot</b>			Full-shot				Baseline	
		MOIRAISmall	MOIRAIBase	MOIRAILarge	PatchTST	TiDE	TFT	DeepAR	AutoARIMA	Seasonal Naive
Electricity	CRPS MSIS	0.072 7.999	0.055 6.172	<u>0.050</u> 5.875	$0.052 {\pm} 0.00 \\ \underline{5.744} {\pm} 0.12$	$0.048 \pm 0.00$ $5.672 \pm 0.08$	$0.050\pm0.00$ $6.278\pm0.24$	$0.065\pm0.01 \\ 6.893\pm0.82$	0.327 29.412	0.070 35.251
Solar	CRPS MSIS	0.471 8.425	<u>0.419</u> <u>7.011</u>	0.406 6.250	0.518±0.09 8.447±1.59	$0.420\pm0.00$ $13.754\pm0.32$	0.446±0.03 8.057±3.51	0.431±0.01 11.181±0.67	1.055 25.849	0.512 48.130
Walmart	CRPS MSIS	0.103 9.371	0.093 8.421	0.098 8.520	$\frac{0.082 \pm 0.01}{6.005 \pm 0.21}$	$0.077 \pm 0.00$ $6.258 \pm 0.12$	$0.087 \pm 0.00$ $8.718 \pm 0.10$	0.121±0.00 12.502±0.03	0.124 9.888	0.151 49.458
Weather	CRPS MSIS	0.049 5.236	<b>0.041</b> 5.136	0.051 <b>4.962</b>	$0.059\pm0.01$ $7.759\pm0.49$	$0.054\pm0.00\ 8.095\pm1.74$	$\frac{0.043 \pm 0.00}{7.791 \pm 0.44}$	0.132±0.11 21.651±17.34	0.252 19.805	0.068 31.293
Istanbul Traffic	CRPS MSIS	0.173 5.937	0.116 4.461	0.112 4.277	0.112±0.00 <b>3.813±0.09</b>	$0.110\pm0.01 \\ 4.752\pm0.17$	$\frac{0.110 \pm 0.01}{4.057 \pm 0.44}$	<b>0.108±0.00</b> 4.094±0.31	0.589 16.317	0.257 45.473
Turkey Power	CRPS MSIS	0.048 7.127	0.040 <u>6.766</u>	0.036 6.341	$0.054\pm0.01 \\ 8.978\pm0.51$	$0.046{\pm}0.01 \\ 8.579{\pm}0.52$	$\frac{0.039 \pm 0.00}{7.943 \pm 0.31}$	$0.066\pm0.02$ $13.520\pm1.17$	0.116 14.863	0.085 36.256

#### 智能决策的新时代

- 企业竞争的核心:数据 + 算法 + 场景
- Transformer让企业从"预测"走向"理解"
- 决策范式转变:
  - 从"经验决策"→"数据驱动决策"→"智能决策"
- 管理者需要理解:
  - 如何在企业中落地小模型(私域知识)
  - 如何治理数据与算法的伦理问题

